

Feature selection \Rightarrow new representation (X)
Feature f ? - original SUBSPACE (not like PCA)

- f correlated with the task = label (classification)
- f correlated with other features?
 - GREEDY - is f correlated with features selected so far?
 - is f redundant w.r.t other features?
- f = spread/variance high?
- f = appropriate for downstream algorithm?

Chi Square test for feature selection

July 23, 2016 Author: david

Feature selection is an important problem in Machine learning. There are many feature selection methods available such as mutual information, information gain, and chi square test. In this post, I will use simple examples to describe how to conduct feature selection using chi square test. I will show that it is easy to use Spark or MapReduce to conduct chi square test based feature selection on large scale data set.

Problem Statement

Suppose there are N instances, and two classes: positive and negative. Given a feature X , we can use Chi Square Test to evaluate its importance to distinguish the class.

By calculating the Chi square scores for all the features, we can rank the features by the chi square scores, then choose the top ranked features for model training.

This method can be easily applied for text mining, where terms or words or N -grams are features. After applying chi square test, we can select the top ranked terms as the features to build a text mining model.

Understand Chi Square Test

Chi Square Test is used in statistics to test the independence of two events.

Given dataset about two events, we can get the observed count o and the expected count E . Chi Square Score measures how much the expected counts E and observed Count o deviate from each other.

In feature selection, the two events are occurrence of the feature and occurrence

∩ per feature
of
indep the values
in a feature column

of the class.

In other words, we want to test whether the occurrence of a specific feature and the occurrence of a specific class are independent.

If the two events are dependent, we can use the occurrence of the feature to predict the occurrence of the class. We aim to select the features, of which the occurrence is highly dependent on the occurrence of the class.

When the two events are independent, the observed count is close to the expected count, thus a small chi square score. So a high value of χ^2 indicates that the hypothesis of independence is incorrect. In other words, the higher value of the χ^2 score, the more likelihood the feature is correlated with the class, thus it should be selected for model training.

How to use Chi Square test for feature selection

Suppose we have a set of training instances that belonging to positive and negative classes. To calculate the χ^2 score of a feature X, we can build the following table, in which there are four numbers:

A: the number of positive instances that contain feature X

B: the number of negative instances that contain feature X

C: the number of positive instances that do not contain feature X

D: the number of negative instances that do not contain feature X

	Label = Y Positive class	Label = N Negative class	total
feature X occurs	A	B	A+B = M
feature X not occurs	C	D	C+D = N - M
total	A+C = P	B+D = N - P	N

Conf table 2x2

pred	class 1	class 2	total
class 1			
class 2			
total			

Count (i,j)

We also define N, M, P, which are describe blow:

N: denote the total number of instances

$M = A + B$: the number of instances that contain feature X

$C + D = N - M$: the number of instances that do not contain feature X

$A + C = P$: the number of positive instances

$B + D = N - P$: the number of negative instances

Let A, B, C, D denotes the observed value, and E_A, E_B, E_C, E_D denote the expected value,

How to calculate the expected value

Based on the null hypothesis that the two events are independent, we can calculate the expected value E_A using the following formula:

$$\rightarrow \frac{E_A}{A+C} = \frac{A+B}{N}$$

So

$$E_A = (A + C) \frac{A+B}{N}$$

The basic idea is that if the two events are independent, the probability that feature X occurs in the Positive class instances should be equal to the probability that feature X occurs in all the instances of the two classes.

Using the similar idea, we can calculate E_B, E_C, E_D .

Using Chi Square Test for feature selection

Using the formula of Chi Square test:

Theory

$$\chi^2 = \frac{1}{d} \sum_{k=1}^n \frac{(O_k - E_k)^2}{E_k}$$

We have

practise

$$\chi^2 = \frac{(A - E_A)^2}{E_A} + \frac{(B - E_B)^2}{E_B} + \frac{(C - E_C)^2}{E_C} + \frac{(D - E_D)^2}{E_D}$$

After some simple calculation we have:

$$\chi^2 = \frac{N(AD - BC)^2}{(A + C)(B + D)(A + B)(C + D)}$$

Given

$$B = M - A$$

$$C = P - A$$

$$D = N - M - (P - A)$$

We have:

same

$$\chi^2 = \frac{N(AN - MP)^2}{PM(N - P)(N - M)}$$

high \Rightarrow select

low \Rightarrow don't select.

How to implement the Chi Square Test Algorithm

HW (A) Original data

\rightarrow Task
Algo (class)

\rightarrow Performance

(B) New feat data
- PCA
- feat select
- feat extraction

\rightarrow Task
Algo/class

\rightarrow Perform

Given a dataset, we can easily obtain:

A: the total number of positive instances that contain feature X,

M: the total number of instances that contain feature X

P: the total number of positive instance,

N: the total number of instance

Apparently, N, P, are constant for all the features,

For each feature, we only need to count A, and M. It is straightforward to implement the algorithm in python or Java for a small dataset.

If we have a large data set with millions of features, we can also easily implement the algorithm using Spark or MapReduce.

Please refer to this article on [using chi square test for feature selection on large scale dataset](#).

A Review of Feature Selection Methods Based on Mutual Information

Jorge R. Vergara · Pablo A. Estévez

Received: date / Accepted: date

Abstract In this work we present a review of the state of the art of information theoretic feature selection methods. The concepts of feature relevance, redundancy and complementarity (synergy) are clearly defined, as well as Markov blanket. The problem of optimal feature selection is defined. A unifying theoretical framework is described, which can retrofit successful heuristic criteria, indicating the approximations made by each method. A number of open problems in the field are presented.

Keywords Feature selection · mutual information · relevance · redundancy · complementarity · synergy · Markov blanket

1 Introduction

Feature selection has been widely investigated and used by the machine learning and data mining community. In this context, a feature, also called attribute or variable, represents a property of a process or system than has been measured, or constructed from the original input variables. The goal of feature selection is to select the smallest feature subset given a certain generalization error, or alternatively finding the best feature subset with k features, that

Jorge R. Vergara

Department of Electrical Engineering, Faculty of Physical and Mathematical Sciences, University of Chile, Chile

Tel.: +56-2-29784207

E-mail: jorgever@ing.uchile.cl

Pablo A. Estévez

Department of Electrical Engineering and Advanced Mining Technology Center, Faculty of Physical and Mathematical Sciences, University of Chile, Chile

Tel.: +56-2-29784207

Fax: +56-2-26953881

E-mail: pestevez@ing.uchile.cl

yields the minimum generalization error. Additional objectives of feature selection are: (i) improve the generalization performance with respect to the model built using the whole set of features, (ii) provide a more robust generalization and a faster response with unseen data, and (iii) achieve a better and simpler understanding of the process that generates the data [31, 23]. We will assume that the feature selection method is used either as a preprocessing step or in conjunction with a learning machine for classification or regression purposes. Feature selection methods are usually classified in three main groups: wrapper, embedded and filter methods [23]. Wrappers [31] use the induction learning algorithm as part of the function evaluating feature subsets. The performance is usually measured in terms of the classification rate obtained on a testing set, i.e., the classifier is used as a black box for assessing feature subsets. Although these techniques may achieve a good generalization, the computational cost of training the classifier a combinatorial number of times becomes prohibitive for high dimensional datasets. In addition, many classifiers are prone to over-learning and show sensitiveness to initialization. Embedded methods [38], incorporate knowledge about the specific structure of the class of functions used by a certain learning machine, e.g. bounds on the leave-one-out error of SVMs [64]. Although usually less computationally expensive than wrappers, embedded methods still are much slower than filter approaches, and the features selected are dependent on the learning machine. Filter methods [17] assume complete independence between the learning machine and the data, and therefore use a metric independent of the induction learning algorithm to assess feature subsets. Filter methods are relatively robust against overfitting, but may fail to select the best feature subset for classification or regression. In the literature, several criteria have been proposed to evaluate single features or feature subsets, among them: inconsistency rate [28], inference correlation [44], classification error [18], fractal dimension [45], distance measure [50, 8], etc. Mutual information (MI) is a measure of statistical independence, that has two main properties. First, it can measure any kind of relation between random variables, including non-linear relationships [14]. Second, MI is invariant under transformations in the feature space that are invertible and differentiable, e.g. translations, rotations and any transformation preserving the order of the original elements of the feature vectors [35, 36]. Many advances in the field have been reported in the last 20 years since the pioneer work of Battiti [4]. Battiti defined the problem of feature selection as the process of selecting the k most relevant variables from an original feature set of m variables, $k < m$. Battiti proposed the greedy selection of a single feature at a time, as an alternative to evaluate the combinatorial explosion of all feature subsets belonging to the original set. The main assumptions of Battiti's work were the following: (a) features are classified as relevant and redundant; (b) an heuristic functional is used to select features, which allows controlling the tradeoff between relevancy and redundancy; c) a greedy search strategy is used; and d) the selected feature subset is assumed optimal. These four assumptions will be revisited in this work to include recent work on a) new definitions on relevant features and other types

of features, b) new information-theoretic functional derived from first principles, c) new search strategies, and d) new definitions of optimal feature subset. In this work, we present a review of filtering feature selection methods based on mutual information, under a unified theoretical framework. We show the evolution of feature selection methods on the last 20 years, describing advantages and drawbacks. The remainder of this work is organized as follows. In section 2 a background on MI is presented. In section 3, the concepts of relevant, redundant and complementary features are defined. In section 4, the problem of optimal feature selection is defined. In section 5, a unified theoretical framework is presented, which allows us to show the evolution of different MI feature selection methods, as well as their advantages and drawbacks. In section 6, a number of open problems in the field are presented. Finally, in section 7, we present the conclusions of this work.

2 Background on MI

→ Intro to Information Theory

2.1 Notation

In this work we will use only discrete random variables, because in practice the variables used in most feature selection problems are either discrete by nature or by quantization. Let F be a feature set and C an output vector representing the classes of a real process. Let's assume that F is the realization of a random sampling of an unknown distribution, where f_i is the i -th variable of F and $f_i(j)$ is the j -th sample of vector f_i . Likewise, c_i is the i -th component of C and $c_i(j)$ is the j -th sample of vector c_i . Uppercase letters denote random sets of variables, and lowercase letters denote individual variables from these sets.

Other notations and terminologies used in this work are the following:

S	Subset of current selected variables.
f_i	Candidate feature to be added to or deleted from the subset of selected features S .
$\{f_i, f_j\}$	Subset composed of the variables f_i and f_j .
$\neg f_i$	All variables in F except f_i . $\neg f_i = F \setminus f_i$.
$\{f_i, S\}$	Subset composed of variable f_i and subset S .
$\neg\{f_i, S\}$	All variables in F except the subset $\{f_i, S\}$. $\neg\{f_i, S\} = F \setminus \{f_i, S\}$
$p(f_i, C)$	Joint mass probability between variables f_i and C .
$ \cdot $	Absolute value / cardinality of a set.

The sets mentioned above are related as follows: $F = f_i \cup S \cup \neg\{f_i, S\}$, $\emptyset = f_i \cap S \cap \neg\{f_i, S\}$. The number of samples in F is n and the total number of variables in F is m .

2.2 Basic Definitions

Entropy, divergence and mutual information are basic concepts defined within information theory [14]. In its origin, information theory was used within the

context of communication theory, to find answers about data compression and transmission rate [52]. Since then, information theory principles have been largely incorporated into machine learning, see for example Principe [47].

2.2.1 Entropy

of distribution

skewness
one val likely

VS
uniformity
all values likely

Entropy (H) is a measure of uncertainty of a random variable. The uncertainty is related to the probability of occurrence of an event. Intuitively, high entropy means that each event has about the same probability of occurrence, while low entropy means that each event has a different probability of occurrence. Formally, the entropy of a discrete random variable x , with mass probability $p(x(i)) = Pr\{x = x(i)\}$, $x(i) \in x$ is defined as:

$$H(x) = - \sum_{i=1}^n p(x(i)) \log_2(p(x(i)))$$

not related to values
(unlike μ, σ)
 $\sum p \cdot \log \frac{1}{p}$

fake die
1 2 3 4 5 6

unif die
1 2 3 4 5 6
 $\frac{1}{6} \frac{1}{6} \frac{1}{6} \frac{1}{6} \frac{1}{6} \frac{1}{6}$
6% 4% 4%

Entropy is interpreted as the expected value of the negative of the logarithm of mass probability. Let x and y be two random discrete variables. The joint entropy of x and y , with joint mass probability $p(x(i), y(j))$, is the sum of the uncertainty contained by the two variables. Formally, joint entropy is defined as follows:

$$H(\{x, y\}) = - \sum_{i=1}^n \sum_{j=1}^n p(x(i), y(j)) \cdot \log_2(p(x(i), y(j)))$$

$p(x,y)$

The joint entropy has values in the range

$$\max(H(x), H(y)) \leq H(\{x, y\}) \leq H(x) + H(y)$$

The maximum value in inequality (3), happens when x and y are completely independent. The minimum value occurs when x is completely dependent on y . The conditional entropy measures the remaining uncertainty of the random variable x when the value of the random variable y is known. The minimum value of the conditional entropy is zero, and it happens when x is statistically dependent on y , i.e., there is no uncertainty in x if we know y . The maximum value happens when x and y are statistically independent, i.e., the variable y does not add information to reduce the uncertainty of x . Formally, the conditional entropy is defined as:

uncertainty of x if we know y

$$H(x|y) = \sum_{j=1}^n p(y(j)) \cdot H(x|y = y(j))$$

where,

$$0 < H(x|y) < H(x)$$

and $H(x|y = y(j))$ is the entropy of all $x(i)$, which are associated with $y = y(j)$.

Another way of representing the conditional entropy is:

$$H(x|y) = H(\{x, y\}) - H(y)$$

$p(x) \cdot p(y) = p(x,y)$ if only if x, y indep

proof $p(x,y) = p(x) \cdot p(y|x) \iff p(y) = p(y|x)$ indep
 $= p(x) \cdot p(y)$

2.2.2 Mutual Information

The mutual information (MI) is a measure of the amount of information that one random variable has about another variable [14]. This definition is useful within the context of feature selection because it gives a way to quantify the relevance of a feature subset with respect to the output vector C . Formally, the MI is defined as follows:

$$I(x; y) = \sum_{i=1}^n \sum_{j=1}^n p(x(i), y(j)) \cdot \log \left(\frac{p(x(i), y(j))}{p(x(i)) \cdot p(y(j))} \right), \quad (7)$$

joint joint
marginals product of marginals

where MI is zero when x and y are statistically independent, i.e., $p(x(i), y(j)) = p(x(i)) \cdot p(y(j))$. The MI is related linearly to entropies of the variables through the following equations:

$$I(x; y) = \begin{cases} H(x) - H(x|y) \\ H(y) - H(y|x) \\ H(x) + H(y) - H(x, y). \end{cases} \quad (8)$$

Fig. 1 shows a Venn diagram with the relationships described in (8).

Let z be a discrete random variable. Its interaction with the other two variables $\{x, y\}$ can be measured by the conditional MI, which is defined as follows:

$$I(x; y|z) = \sum_{i=1}^n p(z(i)) I(x; y|z = z(i)), \quad (9)$$

where $I(x; y|z = z(i))$ is the MI between x and y in the context of $z = z(i)$. The conditional MI allows measuring the information of two variables in the context of a third one, but it does not measure the information among the three variables. Multi-information is an interesting extension of MI, proposed by McGill [42], which allows measuring the interaction among more than two variables. For the case of three variables, the multi-information is defined as follows:

$$I(x; y; z) = \begin{cases} I(\{x, y\}; z) - I(x; z) - I(y; z) \\ I(y; z|x) - I(y; z). \end{cases} \quad (10)$$

The multi-information is symmetrical, i.e., $I(x; y; z) = I(x; z; y) = I(z; y; x) = I(y; x; z) = \dots$. The multi-information has not been widely used in the literature, due to its difficult interpretation, e.g. the multi-information can take negative values, among other reasons. However, there are some interesting papers about the interaction among variables that use this concept [42, 68, 30, 5]. The multi-information can be understood as the amount of information common to all variables (or set of variables), but that is not present in any subset of these variables. To better understand the concept of multi-information within the context of feature selection, let us consider the following example.

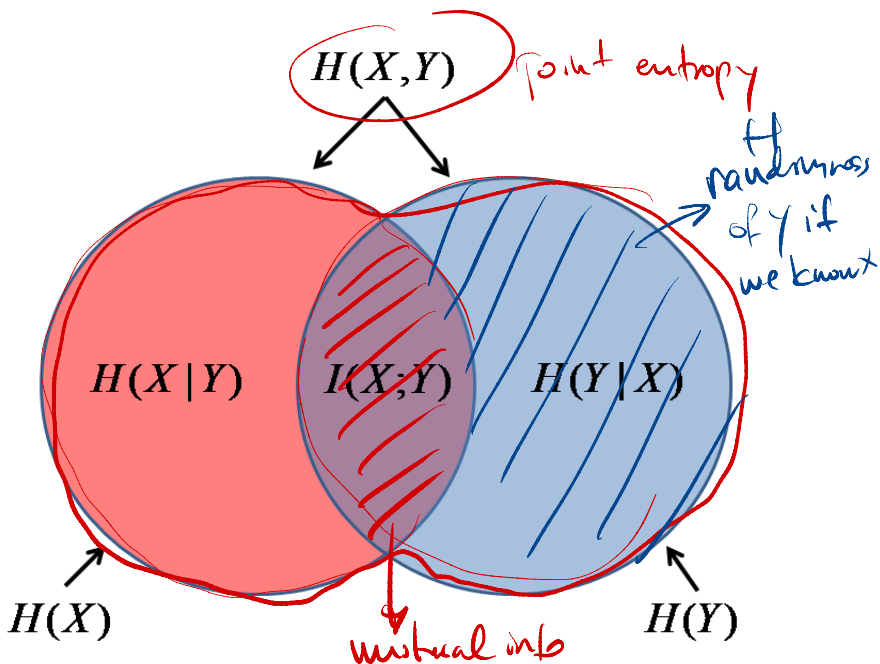


Fig. 1 Venn diagram showing the relations between MI and entropies

Example 1 Let x_1, x_2, x_3 be independent binary random variables. The output of a given system is built through the function $C = x_1 + (x_2 \oplus x_3)$, and $x_4 = x_1$, where $+$ stands for the OR logic function and \oplus represents the XOR logic function.

x_1	x_2	x_3	x_4	$x_2 \oplus x_3$	C
0	0	0	0	0	0
1	0	0	1	0	1
0	1	0	0	1	1
1	1	0	1	1	1
0	0	1	0	1	1
1	0	1	1	1	1
0	1	1	0	0	0
1	1	1	1	0	1

Using eq. (10) to measure the multi-information among x_2, x_3 and C gives: $I(x_2; x_3; C) = I(\{x_2, x_3\}; C) - I(x_2; C) - I(x_3; C)$. Notice that the relevance of single features x_2 and x_3 with respect to C is null, since $I(x_2; C) = I(x_3; C) = 0$, but the joint information of $\{x_2, x_3\}$ with respect to C is greater than zero, $I(\{x_2, x_3\}; C) > 0$. In this case, x_2 and x_3 interact positively to predict C , and

this yields a positive value of the multi-information among these variables. The multi-information among the variables x_1 , x_4 and C is given by: $I(x_1; x_4; C) = I(\{x_1, x_4\}; C) - I(x_1; C) - I(x_4; C)$. The relevance of individual features x_1 and x_4 is the same, i.e., $I(x_1; C) = I(x_4; C) > 0$. In this case the joint information provided by x_1 and x_4 with respect to C is the same as that of each variable acting separately, i.e., $I(\{x_1, x_4\}; C) = I(x_1; C) = I(x_4; C)$. This yields a negative value of the multi-information among these variables. We can deduce that the interaction between x_1 and x_4 does not provide any new information about C . Let us consider now the multi-information among x_1 , x_2 and C , which is zero: $I(x_1; x_2; C) = I(\{x_1, x_2\}; C) - I(x_1; C) - I(x_2; C) = 0$. Since feature x_2 only provides information about C when interacting with x_3 , then $I(\{x_1, x_2\}; C) = I(x_1; C)$. In this case, features x_1 and x_2 do not interact in the knowledge of C .

From the viewpoint of feature selection, the value of the multi-information (positive, negative or zero) gives rich information about the kind of interaction there is among the variables. Let us consider the case where we have a set of already selected features S and a candidate feature f_i , and we measure the multi-information of these variables with the class variable C , $I(f_i; S; C) = I(S; C|f_i) - I(S; C)$. When the multi-information is positive, it means that feature f_i and S are complementary. On the other hand, when the multi-information is negative, it means that by adding f_i we are diminishing the dependence between S and C , because f_i and S are redundant. Finally, when the multi-information is zero, it means that f_i is irrelevant with respect to the dependency between S and C .

The mutual information between a set of m features and the class variable C can be expressed compactly in terms of multi-information as follows:

$$I(\{x_1, x_2, \dots, x_m\}; C) = \sum_{k=1}^m \sum_{\substack{\forall S \subseteq \{x_1, \dots, x_m\} \\ |S| = k}} I([S \cup C]), \quad (11)$$

where $I([S \cup C]) = I(s_1; s_2; \dots; s_k; C)$. Note that the sum on the right side of eq. (11), is taken over all subsets S of size k drawn from the set $\{x_1, \dots, x_m\}$.

3 Relevance, Redundancy and Complementarity

The filter approach to feature selection is based on the idea of relevance, which we will explore in more detail in this section. Basically the problem is to find the feature subset of minimum cardinality that preserves the information contained in the whole set of features with respect to C . This problem is usually solved by finding the relevant features and discarding redundant and irrelevant features. In this section, we review the different definitions of relevance, redundancy and complementarity found in the literature.

3.1 Relevance

Intuitively, a given feature is relevant when either individually or together with other variables, it provides information about C . In the literature there are many definitions of relevance, including different levels of relevance [6, 4, 23, 31, 67, 46, 2, 1, 10, 15]. Kohavi and John [31] used a probabilistic framework to define three levels of relevance: strongly relevant, weakly relevant, and irrelevant features, as shown in Table 1. Strongly relevant features provide unique information about C , i.e., they cannot be replaced by other features. Weakly relevant features provide information about C , but they can be replaced by other features without losing information about C . Irrelevant features do not provide information about C , and they can be discarded without losing information. A drawback of the probabilistic approach is the need of testing the conditional independence for all possible feature subsets, and estimating the probability density functions (pdfs) [48].

An alternative definition of relevance is given under the framework of mutual information [53, 6, 32, 33, 67, 37, 21, 55]. An advantage of this approach is that there are several good methods for estimating MI. The last column of Table 1 shows how the three levels of individual relevance are defined in terms of MI.

Table 1 Levels of relevance for candidate feature f_i , according to probabilistic framework [31] and mutual information framework [43]

Relevance Level	Condition	Probabilistic Approach	Mutual Information Approach
Strongly Relevant	\nexists	$p(C f_i, \neg f_i) \neq p(C \neg f_i)$	$I(f_i; C \neg f_i) > 0$
Weakly Relevant	$\exists S \subset \neg f_i$	$p(C f_i, \neg f_i) = p(C \neg f_i)$ \wedge $p(C f_i, S) \neq p(C S)$	$I(f_i; C \neg f_i) = 0$ \wedge $I(f_i; C S) > 0$
Irrelevant	$\forall S \subseteq \neg f_i$	$p(C f_i, S) = p(C S)$	$I(f_i; C S) = 0$

The definitions shown in Table 1 give rise to several drawbacks, which are summarized as follows:

1. To classify a given feature f_i , as irrelevant, it is necessary to assess all possible subsets S of $\neg f_i$. Therefore this procedure is subject to the curse of dimensionality [7, 57].
2. The definition of strongly relevant features is too restrictive. If two features provides information about the class but are redundant, then both features will be discarded by this criterion. For example, let $\{x_1, x_2, x_3\}$ be a set of 3 variables, where $x_1 = x_2$, and x_3 is noise, and the output class is defined as $C = x_1$. Following the strong relevance criterion we have $I(x_1; C|\{x_2, x_3\}) = I(x_2; C|\{x_1, x_3\}) = I(x_3; C|\{x_1, x_2\}) = 0$.

- The definition of weak relevance is not enough for deciding whether to discard a feature from the optimal feature set. It is necessary to discriminate between redundant and non-redundant features.

3.2 Redundancy

Yu and Liu [67] proposed a finer classification of features into weakly relevant but redundant and weakly relevant but non-redundant. Moreover, the authors defined the set of optimal features as the one composed by strongly relevant features and weakly relevant but non-redundant features. The concept of redundancy is associated with the level of dependency among two or more features. In principle we can measure the dependency of a given feature f_i with respect to a feature subset $S \subseteq \neg f_i$, by simply using the MI, $I(f_i; S)$. This information theoretic measure of redundancy satisfies the following properties: it is symmetric, non-linear, non-negative, and does not diminish when adding new features [43]. However, using this measure it is not possible to determine concretely with which features of S is f_i redundant. This calls for more elaborated criteria of redundancy, such as the Markov blanket [33, 67], and total correlation [62]. The Markov blanket is a strong condition for conditional independence, and is defined as follows.

Definition 1 (Markov blanket) Given a feature f_i , the subset $M \subseteq \neg f_i$ is a Markov blanket of f_i iff [33, 67]:

$$p(\{F \setminus \{f_i, M\}, C\} | \{f_i, M\}) = p(\{F \setminus \{f_i, M\}, C\} | M). \quad (12)$$

This condition requires that M subsumes all the information that f_i has about C , but also about all other features $\{F \setminus \{f_i, M\}\}$. It can be proved that strongly relevant features do not have a Markov blanket [67].

The Markov blanket condition given by Eq. (12) can be rewritten in the context of information theory as follows [43]:

$$I(f_i; \{C, \neg f_i, M\} | M) = 0. \quad (13)$$

An alternative measure of redundancy is the total correlation or multivariate correlation [62]. Given a set of features $F = \{f_1, \dots, f_m\}$, the total correlation is defined as follows:

$$C(f_1; \dots; f_m) = \sum_{i=1}^m H(f_i) - H(f_1, \dots, f_m). \quad (14)$$

Total correlation measures the common information (redundancy) among all the variables in F . If we want to measure the redundancy between a given variable f_i and any feature subset $S \subseteq \neg f_i$, then we can use the total correlation as:

$$C(f_i; S) = H(f_i) + H(S) - H(f_i, S), \quad (15)$$

however this corresponds to the classic definition of MI, i.e., $C(f_i; S) = I(f_i; S)$.

3.3 Complementarity

The concept of complementarity has been re-discovered several times [43, 9, 10, 61, 12]. Recently, it has become more relevant because of the development of more efficient techniques to estimate MI in high-dimensional spaces [34, 27]. Complementarity, also known as synergy, measures the degree of interaction between an individual feature f_i and feature subset S given C , through the following expression ($I(f_i; S|C)$). To illustrate the concept of complementarity, we will start expanding the multi-information among f_i , C and S . Decomposing the multi-information in its three possible expressions we have:

$$I(f_i; S; C) = \begin{cases} I(f_i; S|C) - I(f_i; S) \\ I(f_i; C|S) - I(f_i; C) \\ I(S; C|f_i) - I(S; C). \end{cases} \quad (16)$$

According to eq. (16), the first row shows that the multi-information can be expressed as the difference between complementarity ($I(f_i; S|C)$) and redundancy ($I(f_i; S)$). A positive value of the multi-information entails a dominance of complementarity over redundancy. Analyzing the second row of eq. (16), we observe that this expression becomes positive when the information that f_i has about C is greater when it interacts with subset S with respect to the case when it does not. This effect is called complementarity. The third row of eq. (16), gives us another viewpoint of the complementarity effect. The multi-information is positive when the information that S has about C is greater when it interacts with feature f_i compared to the case when it does not interact. Assuming that the complementarity effect is dominant over redundancy, Fig. 2 illustrates a Venn diagram with the relationships among complementarity, redundancy and relevancy.

4 Optimal Feature Subset

In this section we review the different definitions of the optimal feature subset, S_{opt} , given in the literature, as well as the search strategies used for obtaining this optimal set. According to [58], in practice the feature selection problem must include a classifier or an ensemble of classifiers, and a performance metric. The optimal feature subset is defined as the one that maximizes the performance metric having minimum cardinality. However, filter methods are independent of both the learning machine and the performance metric. Any filter method corresponds to a definition of relevance that employs only the data distribution [58]. Yu and Liu [67] defined the optimal feature set as composed of all strongly relevant features and the weakly relevant but not redundant features. In this section we review the definitions of the optimal feature subset from the viewpoint of filter methods, in particular MI feature selection methods. The key notion is conditional independence, which allows defining the sufficient feature subset as follows [6, 24]:

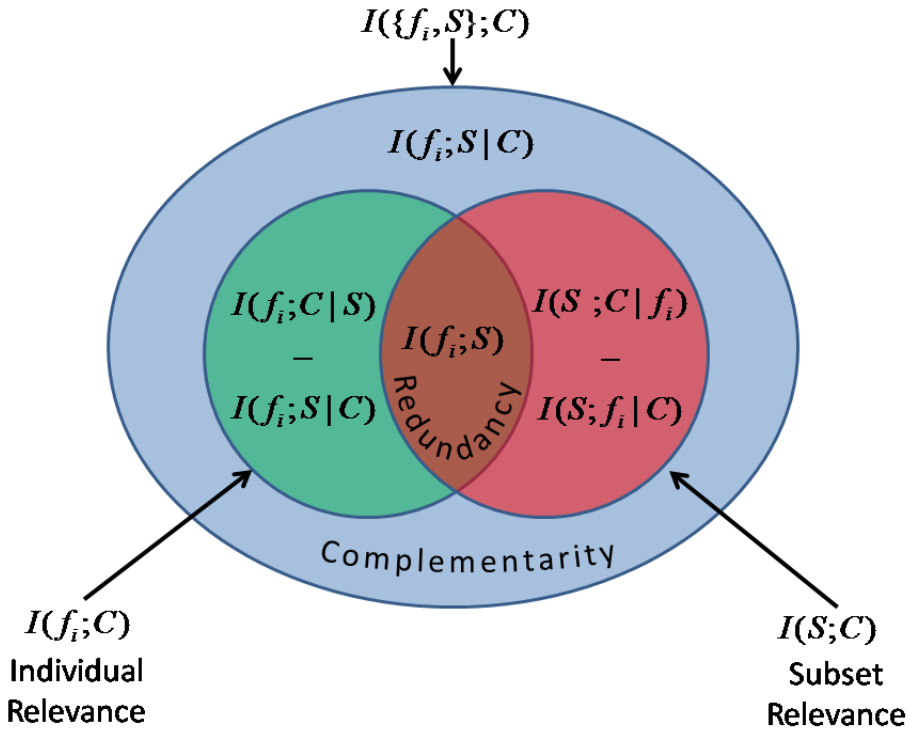


Fig. 2 Venn diagram showing the relationships among complementarity, redundancy and relevancy, assuming that the multi-information among f_i , S and C is positive.

Definition 2 $S \subseteq F$ is a sufficient feature subset iff

$$p(C|F) = p(C|S). \quad (17)$$

This definition implies that C and $\neg S$ are conditionally independent, i.e., $\neg S$ provides no additional information about C in the context of S . However, we still need a search strategy to select the feature subset S , and an exhaustive search using this criterion is impractical due to the curse of dimensionality.

In probability the measure of sufficient feature subset can be expressed as the expected value over $p(F)$ of the Kullback-Leibler divergence between $p(C|F)$ and $p(C|S)$ [33]. According to Guyon *et al.* [24], this can be expressed in terms of MI as follows:

$$DMI(S) = I(F; C) - I(S; C). \quad (18)$$

Guyon *et al.* [24] proposed solving the following optimization problem:

$$\min_{S \subseteq F} |S| + \lambda \cdot DMI(S), \quad (19)$$

where $\lambda > 0$ represents the Lagrange multiplier. If S is a sufficient feature subset, then $DMI(S) = 0$, and eq. (19) is reduced to $\min_{S \subseteq F} |S|$. Since $I(F; C)$ is constant, eq. (19) is equivalent to:

$$\min_{S \subseteq F} |S| - \lambda \cdot I(S; C). \quad (20)$$

The feature selection problem corresponds to finding the smallest feature subset that maximizes $I(S; C)$. Since the term $\min_{S \subseteq F} |S|$ is discrete, the optimization of (20) is difficult. Tishby *et al.* [55] proposed replacing the term $\min_{S \subseteq F} |S|$ with $I(F; S)$.

An alternative approach to optimal feature subset selection is using the concept of the Markov blanket (MB). Remember that the Markov blanket, M , of a target variable C , is the smallest subset of F such that C is independent of the rest of the variables $F \setminus M$. Koller and Sahami [33] proposed using MBs as the basis for feature elimination. They proved that features eliminated sequentially based on this criterion remain unnecessary. However, the time needed for inducing an MB grows exponentially with the size of this set, when considering full dependencies. Therefore most MB algorithms implement approximations based on heuristics, e.g. finding the set of k features that are strongly correlated with a given feature [33]. Fast MB discovery algorithms have been developed for the case of distributions that are faithful to a Bayesian Network [58, 59]. However, these algorithms require that the optimal feature subset does not contain multivariate associations among variables, which are individually irrelevant but become relevant in the context of others [11]. In practice, this means for example that current MB discovery algorithms cannot solve Example 1 due to the XOR function.

An important caveat is that both feature selection approaches, sufficient feature subset and MBs, are based on estimating the probability distribution of C given the data. Estimating posterior probabilities is a harder problem than classification, e.g. in using a 0/1-loss function only the most probable classification is needed. Therefore, this effect may render some features contained in sufficient feature subset or in the MB of C unnecessary [58, 56, 24].

4.1 Relation between MI and Bayes error classification

There are some interesting results relating the MI between a random discrete variable f and a random discrete target variable C , with the minimum error obtained by maximum a posteriori classifier (Bayer classification error) [26, 14, 20]. The Bayes error is bounded above and below according to the following expression:

$$1 - \frac{I(f; C) + \log(2)}{\log(|C|)} \leq e_{bayes}(f) \leq \frac{1}{2} (H(C) - I(f; C)). \quad (21)$$

Interestingly, Eq. (21) shows that both limits are minimized when the MI, $I(f; C)$, is maximized.

4.2 Search strategies

According to Guyon *et al.* [24], a feature selection method has three components: 1) Evaluation criterion definition, e.g. relevance for filter methods, 2) evaluation criterion estimation, e.g. sufficient feature selection or MB for filter methods, and 3) search strategies for feature subset generation. In this section, we briefly review the main search strategies used by MI feature selection methods. Given a feature set F of cardinality m , there are 2^m possible subsets, therefore an exhaustive search is impractical for high-dimensional datasets.

There are two basic search strategies: optimal methods and sub-optimal methods [63]. Optimal search strategies include exhaustive search and accelerated methods based on the monotonic property of a feature selection criterion, such as branch and bound. But optimal methods are impractical for high-dimensional datasets, therefore sub-optimal strategies must be used.

Most popular search methods are **sequential forward selection** (SFS) [65] and sequential backward elimination (SBE) [41]. Sequential forward selection is a bottom-up search, which starts with an empty set, and adds new features one at a time. Formally, it adds the candidate feature f_i that maximizes $I(S; C)$ to the subset of selected features S , i.e.,

$$S = S \cup \left\{ \arg \max_{f_i \in F \setminus S} (I(\{S, f_i\}; C)) \right\}. \quad (22)$$

→ already selected
→ label
→ new feat candidate

Greedy build up

Sequential backward elimination is a top-down approach, which starts with the whole set of features, and deletes one feature at a time. Formally, it starts with $S = F$, and proceeds deleting the less informative features one at a time, i.e.,

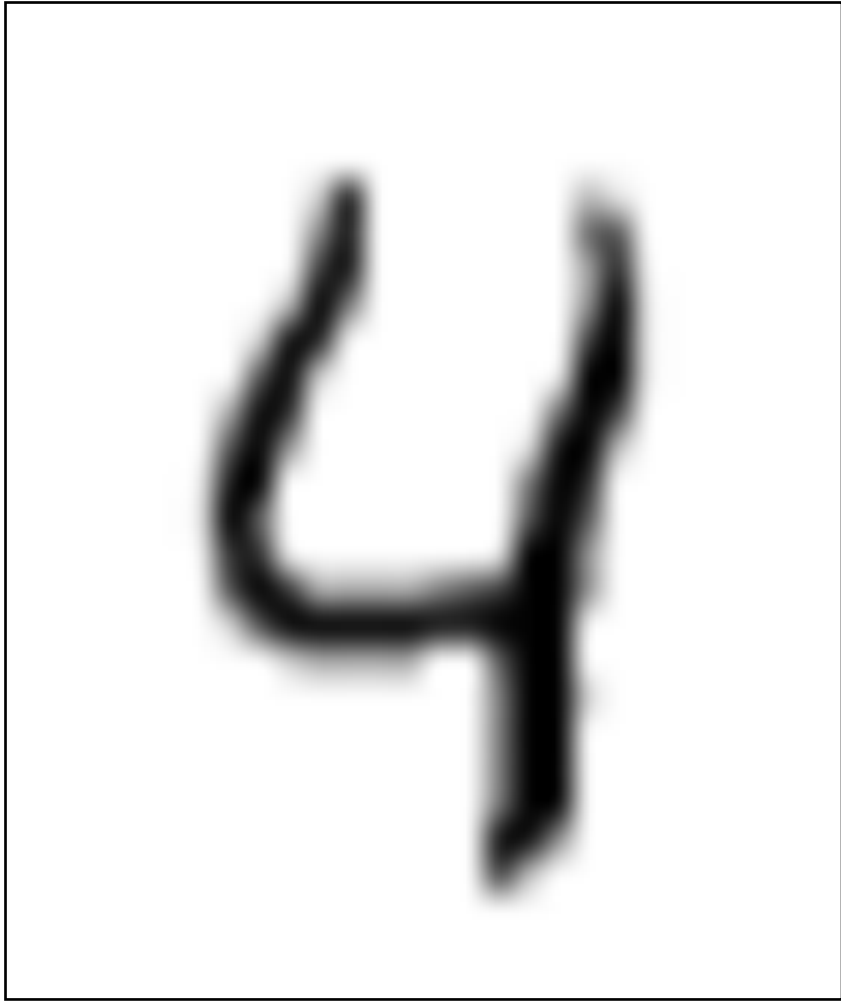
$$S = S \setminus \left\{ \arg \min_{f_i \in S} (I(S \setminus f_i; C)) \right\}. \quad (23)$$

→ Greedy elimination

Usually **backward elimination is computationally more expensive** than forward selection, e.g. when searching for a small subset of features. However, backward elimination can usually find better feature subsets, because most forward selection methods do not take into account the relevance of variables in the context of features not yet included in the subset of selected features [23]. Both kinds of searching methods suffer from the nested effect, meaning that in forward selection a variable cannot be deleted from the feature set once it has been added, and in backward selection a variable cannot be re-incorporated once it has been deleted. Instead of adding a single feature at a time, some generalized forward selection variants add several features, to take into account the statistical relationship between variables [63]. Likewise, the generalized backward elimination deletes several variables at a time. An enhancement may be obtained by combining forward and backward selection, avoiding the nested effect. The strategy “plus-1-take-away- r ” [54] adds to S l features and then removes the worst r features if $l > r$, or deletes r features and then adds l features if $r < l$.

**HAAAR-like features for
images**

Images

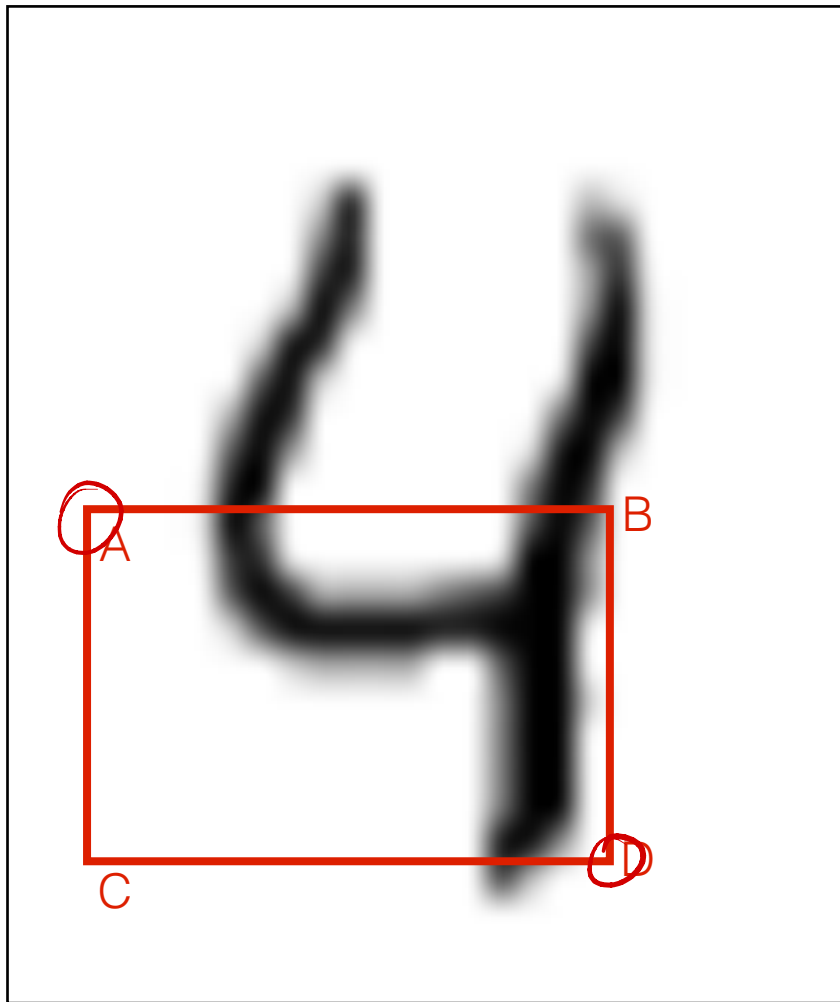


- digit images are scanned hand written digits

Digit scan dataset

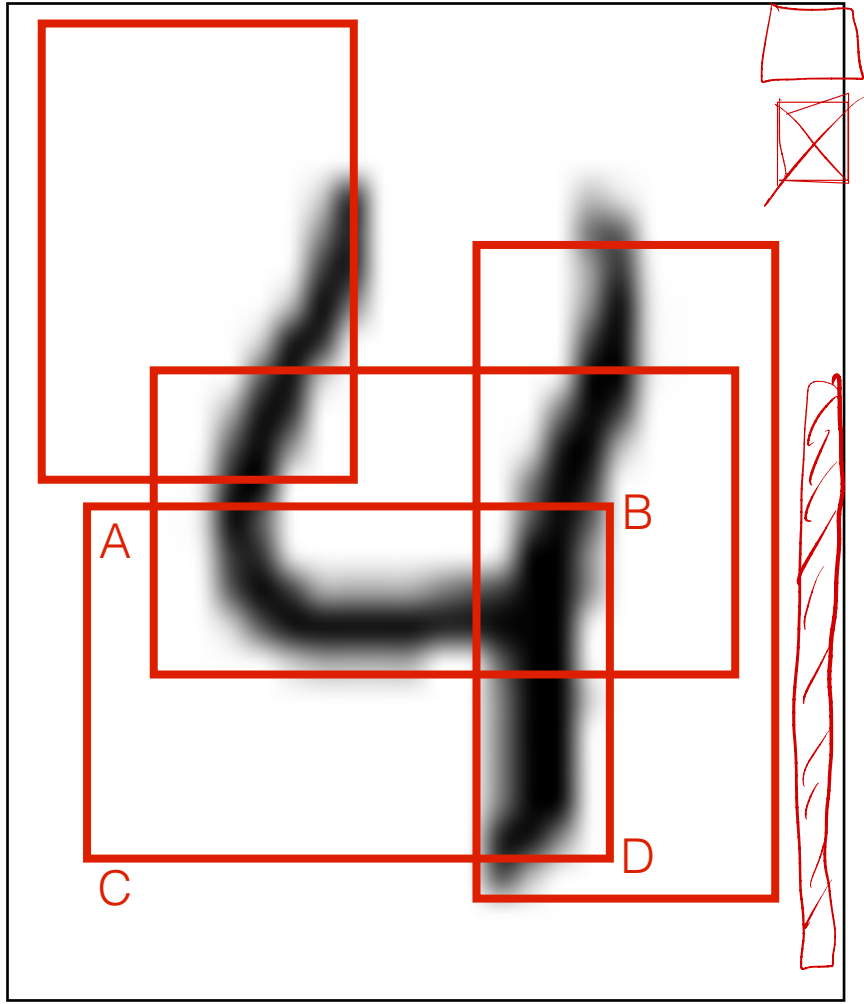
- 60,000 scans
- 10 classes : 0,1,2,...,9
 - roughly uniform distributed
- each scanned image 28x28 pixels square
- comes split into (train, test)
 - no cross validation
- very learnable: most algorithms score 5% or less error
- <http://yann.lecun.com/exdb/mnist/>

Rectangle black level



- rectangle ABCD can act like an image “mask” : it selects/cuts that rectangle out of an image
 - or of any image
- stat ex total amount of black*

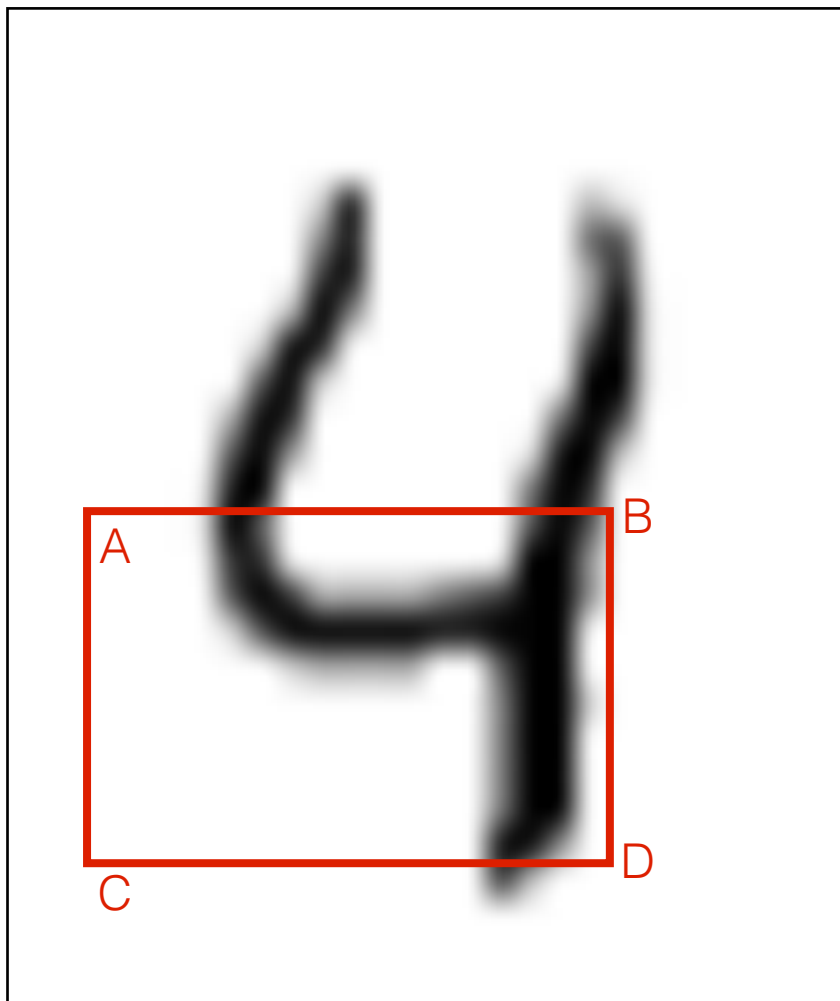
Rectangle black level



- a given set S of rectangles cuts S different masks for an image

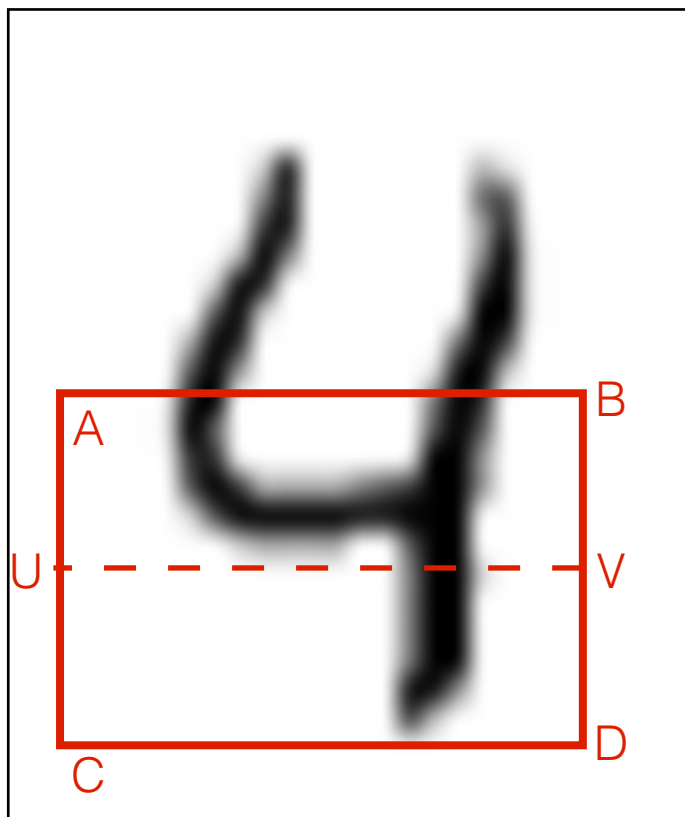
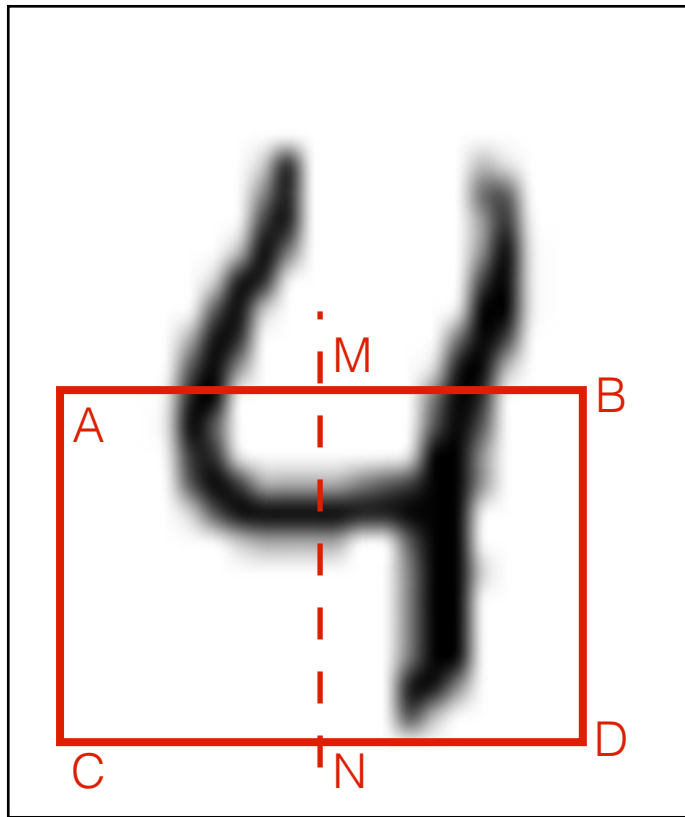
Rectangle black level

0



- for each rectangle $r=ABCD$ on image X we can compute a “black value”
 - $\text{black}_r(X)$ = number of black pixels in the mask cut by r in image X
- we can compute $\text{black}_r(X)$ efficiently, if we compute in the right order!
 - dynamic programming

Vertical, horizontal features for a rectangle



- horizontal feature

$$\begin{aligned}\Delta_{hr}(X) &= \text{black}_{r\text{-left}}() - \text{black}_{r\text{-right}}(X) \\ &= \text{black}_{AMCN}(X) - \text{black}_{MBND}(X)\end{aligned}$$

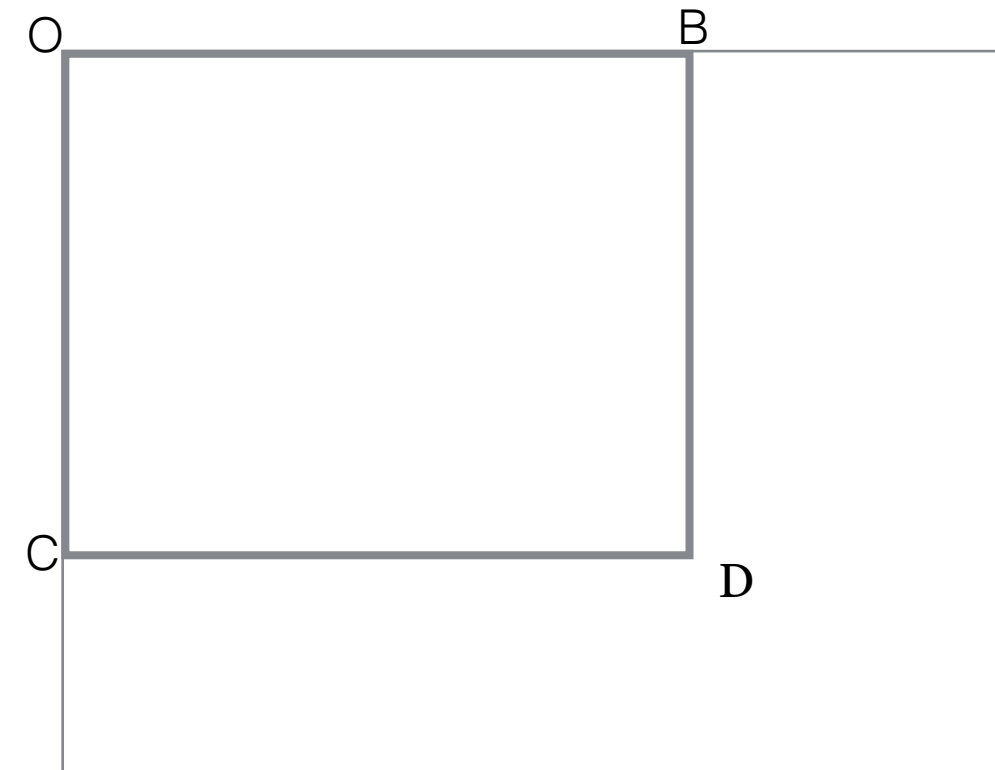
- vertical feature

$$\begin{aligned}\Delta_{hv}(X) &= \text{black}_{r\text{-top}}() - \text{black}_{r\text{-bottom}}(X) \\ &= \text{black}_{ABUV}(X) - \text{black}_{UVCD}(X)\end{aligned}$$

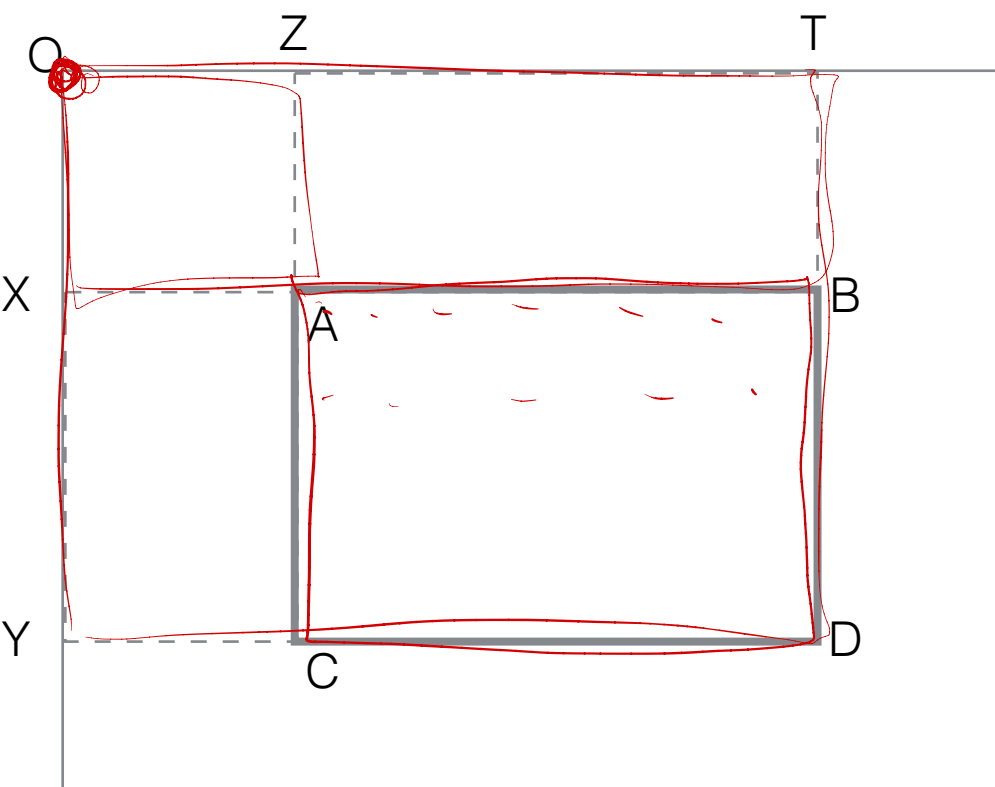
- $|S|$ rectangles, 2 features each $\Rightarrow 2|S|$ features extracted (from each image)

- if we also store the $\text{black}_r(X)$ value, that's 3 features/rectangle ($\text{black}_r(X)$, $\Delta_{hr}(X)$, $\Delta_{hv}(X)$) for $3|S|$ features extracted.

How to compute $\text{black}_r(X)$ efficiently

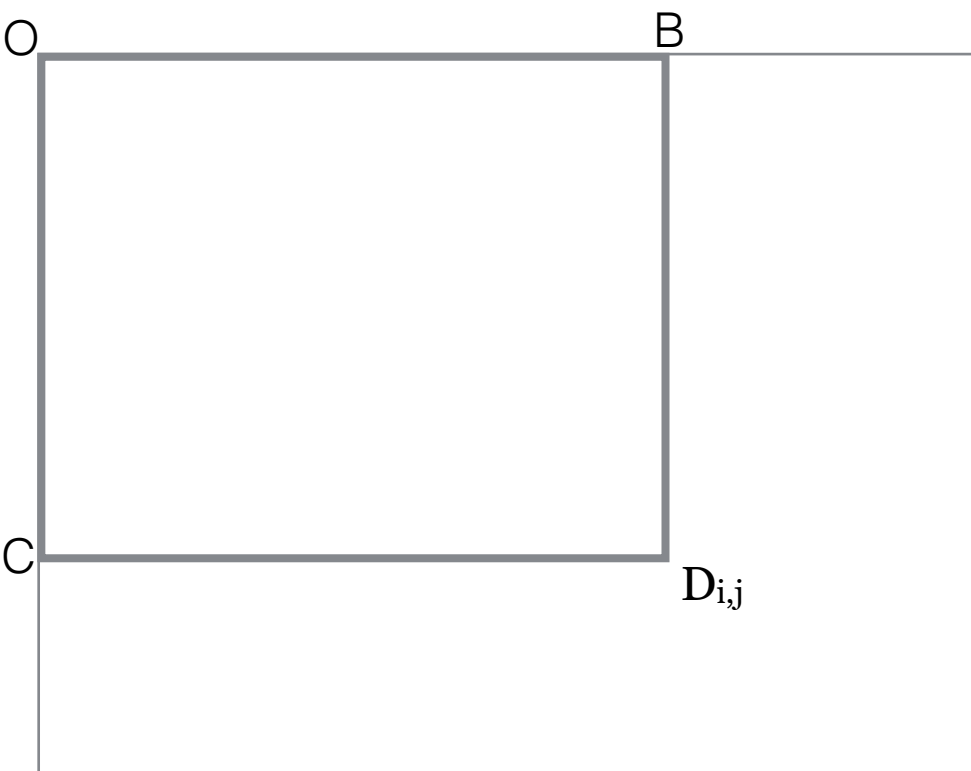


- first compute it for all rectangles cornered in O (A=O) fix image corner.
 - That is compute $\text{black}_r(X)$ for each pixel D



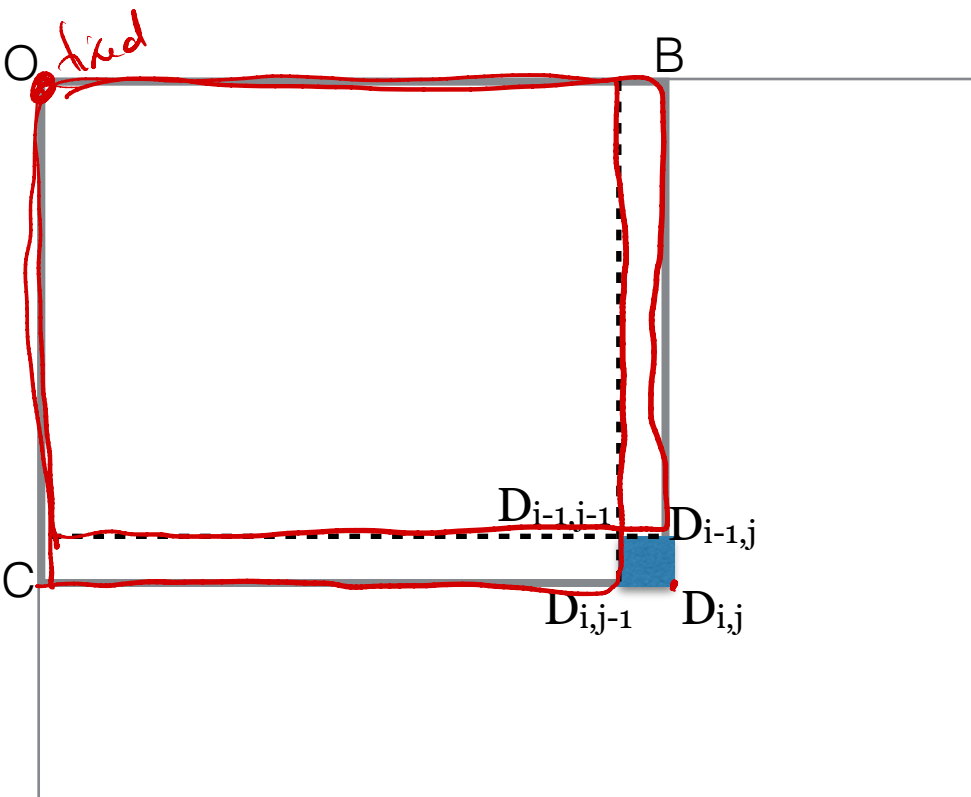
- then every rectangle $r=ABCD$ can be computed in constant time from O-cornered rectangles
- $\text{black}(\text{rectangle } ABCD) = \text{black}(OTYD) - \text{black}(OTXB) - \text{black}(OZYC) + \text{black}(OZXA)$

O-corner rectangles computation



- $r=OBCD$ determined by D
- naively one can compute all $\text{black}_r(X) = \text{black}_D(X)$ for all rectangles as
 - for $i=1:n$
 - for $j=1:n$
 - $D=D_{ij}$ pixel
 - $\text{black}_{D_{ij}}(X) = \text{count of black pixels in } OBCD$
- total $O(n^4)$ running time
 - $n = \text{size of the square image}$

O-corner rectangles : dynamic programming



- $r=OBCD$ determined by D
- dynamic programming computes a rectangle from the rectangle computed already

- for $i=1:n$

- for $j=1:n$

- $D=D_{ij}$ pixel

$$\begin{aligned} \text{black_}D_{ij}(X) &= \\ &\text{black_}D_{i,j-1}(X) + \\ &\text{black_}D_{i-1,j}(X) - \\ &\text{black_}D_{i-1,j-1}(X) + \\ &\text{black}(\text{pixel_}D_{ij}, X) \end{aligned}$$

- total $O(n^2)$ running time
 - much better

5 A Unified Framework for Mutual Information Feature Selection

Many MI feature selection methods have been proposed in the last 20 years. Most methods define heuristic functionals to assess feature subsets combining definitions of relevant and redundant features. Brown *et al.* [10] proposed a unifying framework for information theoretic feature selection methods. The authors posed the feature selection problem as a conditional likelihood of the class labels, given features. Under the filter assumption [10], conditional likelihood is equivalent to conditional mutual information (CMI), i.e., the feature selection problem can be posed as follows:

$$\begin{aligned} \min_{S \subseteq F} |S| & \quad (24) \\ \text{subject to: } \min_{S \subseteq F} I(\neg S; C|S). \end{aligned}$$

This corresponds to the smallest feature subset such that the CMI is minimal. Starting from this objective function, the authors used MI properties to deduce some common heuristic criteria used for MI feature selection. Several criteria can be unified under the proposed framework. In particular, they showed that common heuristics based on linear combinations of information terms, such as Battiti’s MIFS [4], conditional infomax feature extraction (CIFE) [40, 22], minimum-redundancy maximum relevance (mRMR) [46], and joint mutual information (JMI) [66], are all low-order approximations to the conditional likelihood optimization problem. However, the unifying framework proposed by Brown *et al.* [10] fell short of deriving (explaining) non-linear criteria using min or max operators such as Conditional Mutual Information Maximization (CMIM) [21], Informative Fragments [61], and ICAP [29].

Let us start with the assumption that $I(F; C)$ measures all the information about the target variable contained in the set of features. This assumption is based on the additivity property of MI [14, 32], which states that the information about a given system is maximal when all features (F) are used to estimate the target variable (C). Using the chain rule, $I(F; C)$ can be decomposed as follows:

$$I(F; C) = I(S; C) + I(\neg S; C|S). \quad (25)$$

As $I(F; C)$ is constant, maximizing $I(S; C)$ is equivalent to minimizing $I(\neg S; C|S)$. Many MI feature selection methods maximize the first term on the right side of (25). This is known as the criterion of maximal dependency (MD) [46]. On the other hand, other criteria are based on the idea of minimizing the CMI, i.e. the second term on the right hand side of eq. (25).

In the following we describe the approach of Brown *et al.* [10] for deriving sequential forward selection and sequential backward elimination algorithms, which are based on minimizing the CMI. For the convenience of the reader, we present the equivalent procedure in parallel when maximizing dependency (MD). In practice, a search strategy is needed to find the best feature subset. As we saw in section 4.2 the most popular methods are sequential forward selection and sequential backward elimination. Before proceeding we need to define some notation.

Table 2 Parallel between MD and CMI approaches for sequential forward selection

MD	CMI
$\max_{f_i \in \neg S^t} I(S^{t+1}; C) =$	$\min_{f_i \in \neg S^t} I(\neg S^{t+1}; C S^{t+1})$
$\max_{f_i \in \neg S^t} I(\{S^t, f_i\}; C) =$	$\min_{f_i \in \neg S^t} I(\neg S^t \setminus f_i; C \{S^t, f_i\})$
$\max_{f_i \in \neg S^t} I(S^t; C) \overset{a}{\underline{+}} \max_{f_i \in \neg S^t} I(f_i; C S^t)$	$\min_{f_i \in \neg S^t} I(\neg S^t; C S^t) \overset{b}{\underline{+}} \min_{f_i \in \neg S^t} (-I(f_i; C S^t))$
\Downarrow	\Downarrow
$\max_{f_i \in \neg S^t} I(f_i; C S^t)$	$\max_{f_i \in \neg S^t} I(f_i; C S^t)$

^a This term is independent of f_i .

^b This term has the same value $\forall f_i$.

S^t	Subset of selected variables at time t .
f_i	Candidate feature to be added to or eliminated from feature subset S^t at time t .
	$f_i = \arg \max_{f_i \in \neg S^t} I(f_i; C S^t)$ in forward selection.
	$f_i = \arg \min_{f_i \in S^t} I(f_i; C S^t \setminus f_i)$ in backward elimination.
s_j	A given feature in S^t .
$\neg s_j$	The complement set of feature s_j with set S^t , i.e., $\neg s_j = S^t \setminus s_j$
S^{t+1}	Subset of selected variables at time $t+1$.
	$S^{t+1} \leftarrow \{S^t, f_i\}$ in forward selection.
	$S^{t+1} \leftarrow S^t \setminus f_i$ in backward elimination.
$\neg S^{t+1}$	Complement of feature subset S^{t+1} , i.e. $F = \{S^{t+1}, \neg S^{t+1}\}$.
	$\neg S^{t+1} \leftarrow \{\neg S^t \setminus f_i\}$ in forward selection.
	$\neg S^{t+1} \leftarrow \{\neg S^t, f_i\}$ in backward elimination.

Table 2 shows that for the case of sequential forward selection, we achieve the same result when using the MD or CMI approach: the SFS algorithm consists of maximizing $I(f_i; C | S^t)$. Analogously, Table 3 shows that for the case of sequential backward elimination, again we achieve the same result when using MD or CMI approaches: the SBE algorithm consists of minimizing $I(f_i; C | S^t \setminus f_i)$.

For space limitations, we will develop here only the case of forward feature selection, but the procedure is analogous for the case of backward feature elimination. The expression $I(f_i; C | S^t)$ can be expanded as follows [12]:

$$I(f_i; C | S^t) = I(f_i; C) - I(f_i; S^t) + I(f_i; S^t | C). \quad (26)$$

The first term on the right hand side of (26) measures the individual relevance of the candidate feature f_i with respect to output C ; the second term measures the redundancy of the candidate feature with the feature subset of previously selected features S^t ; and the third term measures the complementarity between S^t and f_i in the context of C . However, from the practical point of view, eq. (26) presents the difficulty of estimating MI in high-dimensional spaces, due to the presence of the set S^t in the second and third terms.

Table 3 Parallel between MD and CMI approaches for sequential backward elimination

MD	CMI
$\max_{f_i \in S^t} I(S^{t+1}; C) =$	$\min_{f_i \in S^t} I(\neg S^{t+1}; C S^{t+1})$
$\max_{f_i \in S^t} I(S^t \setminus f_i; C) =$	$\min_{f_i \in S^t} I(\{\neg S^t, f_i\} \setminus f_i; C S^t \setminus f_i)$
$\max_{f_i \in S^t} I(S^t; C) \stackrel{a}{=} + \max_{f_i \in S^t} (-I(f_i; C S^t \setminus f_i))$	$\min_{f_i \in S^t} I(\neg S^t; C S^t) \stackrel{b}{=} + \min_{f_i \in S^t} (I(f_i; C S^t \setminus f_i))$
\Downarrow	\Downarrow
$\min_{f_i \in S^t} I(f_i; C S^t \setminus f_i)$	$\min_{f_i \in S^t} I(f_i; C S^t \setminus f_i)$

^a This term is independent of f_i .

^b This term has the same value $\forall f_i$.

In what follows, we take a detour from the derivation of Brown *et al.* [10], using our own alternative approach. To avoid the previously mentioned problem, $I(f_i; S^t)$ with $|S^t| = p$ can be calculated by averaging all expansions over every single feature in S , by using the chain rule as follows:

$$\begin{aligned}
 I(f_i; S^t) &= I(f_i; s_1) + I(f_i; \neg s_1 | s_1) \\
 I(f_i; S^t) &= I(f_i; s_2) + I(f_i; \neg s_2 | s_2) \\
 \vdots &= \vdots \\
 I(f_i; S^t) &= I(f_i; s_p) + I(f_i; \neg s_p | s_p)
 \end{aligned}$$

$$I(f_i; S^t) = \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; s_j) + \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; \neg s_j | s_j). \quad (27)$$

Analogously, we can obtain the following expansion for the conditional mutual information, $I(f_i; S^t | C)$:

$$I(f_i; S^t | C) = \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; s_j | C) + \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; \neg s_j | \{C, s_j\}). \quad (28)$$

Substituting (27) and (28) into eq. (26) yields:

$$\begin{aligned}
 I(f_i; C | S^t) &= I(f_i; C) - \left(\frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; s_j) + \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; \neg s_j | s_j) \right) \\
 &\quad + \left(\frac{1}{|S|} \sum_{s_j \in S} I(f_i; s_j | C) + \frac{1}{|S|} \sum_{s_j \in S} I(f_i; \neg s_j | \{C, s_j\}) \right).
 \end{aligned} \quad (29)$$

Eq. (29), can be approximated by considering assumptions of lower-order dependencies between features [3]. Features $s_j \in S^t$ are assumed to have only

one-to-one dependencies with f_i or C . Formally, assuming statistical independence:

$$\begin{aligned} p(f_i|S^t) &= \prod_{s_j \in S^t} p(f_i|s_j) \\ p(f_i|\{S^t, C\}) &= \prod_{s_j \in S^t} p(f_i|\{s_j, C\}), \end{aligned} \quad (30)$$

we obtain the following low-order approximation:

$$I(f_i; C|S^t) \approx I(f_i; C) - \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; s_j) + \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; s_j|C). \quad (31)$$

Notice that eq. (31) is an approximation of the multidimensional MI expressed by eq. (26). Interestingly, Brown *et al.* (10) deduced a similar formula but with coefficients $1/|S^t|$ replaced by unity constants.

Eq. (31) allows deriving some well-known heuristic feature selection methods. When only the first two terms of Eq. (31) are taken into account, it corresponds exactly to the minimal redundancy maximal relevance (mRMR) criterion proposed in (46). Moreover, if the term $1/|S^t|$ is replaced by a user defined parameter β , then we obtain the MIFS criterion (*Mutual Information Feature Selection*) proposed by Battiti (4). When considering only the first term in eq. (31), we obtain the MIM criterion (39).

Eq. (31) with its three terms corresponds exactly to the Joint Mutual Information (JMI) (66, 10). Also it corresponds with the Conditional Infomax Feature Extraction (CIFE) criterion proposed in (40), when the coefficient $|S^t| = 1, \forall t$. Moreover, the Conditional Mutual Information based Feature Selection (CMIFS) criterion proposed in (12) is an approximation of eq. (29), where only 0, 1 or 2 out of t summation terms are considered in each term. The CMIFS criterion is the following:

$$J_{cmifs}(f_i) = I(f_i; C) - I(f_i; s_t) + \sum_{s_j \in S; j \in \{1, t\}} I(f_i; s_j|C) - I(f_i; s_t|s_1). \quad (32)$$

The previously mentioned methods do not take into account the terms containing $\neg s_j$ in eq. (29). This entails the assumption that f_i and $\neg s_j$ are independent, therefore ($I(f_i; \neg s_j) = I(f_i; \neg s_j|C) = 0$). This approximation can generate errors in the sequential selection or backward elimination of variables. In order to somehow take into account the missing terms, let us consider the following alternative approximation of $I(f_i; C|S^t)$:

$$\begin{aligned} I(f_i; C|S^t) &= I(f_i; C) + I(f_i; S^t; C) = \\ &= I(f_i; C) + I(f_i; \{s_j, \neg s_j\}; C) = \\ &= I(f_i; C) + I(f_i; s_j; C) + I(f_i; \neg s_j; C|s_j) = \\ &= I(f_i; C|s_j) + I(f_i; \neg s_j; C|s_j). \end{aligned} \quad (33)$$

Averaging this decomposition over every single feature $s_j \in S^t$ we have:

$$I(f_i; C|S^t) = \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; C|s_j) + \frac{1}{|S^t|} \sum_{s_j \in S^t} I(f_i; \neg s_j; C|s_j). \quad (34)$$

The Interaction Capping (ICAP) [29] criterion approximates eq. (33) by the following expression:

$$J_{icap}(f_i) = I(f_i; C) + \sum_{s_j \in S} \min(0, I(f_i; s_j; C)). \quad (35)$$

In ICAP [29], the information of variable f_i is penalized when the interaction between f_i , s_j and C becomes redundant ($I(f_i; s_j; C) < 0$), but the complementarity relationship among variables is neglected when $I(f_i; s_j; C) > 0$. The authors considered a Naive Bayes classifier, which assumes independence between variables.

Eq. (34) allows deriving the Conditional Mutual Information Maximization (CMIM) criterion [21], when we consider only the first term on the right hand side of this equation and replace the mean operator with a minimum operator. CMIM discards the second term on the right hand side of eq. (34) completely, taking into account only one-to-one relationships among variables and neglecting the multi-information among f_i , $\neg s_j$ and C in the context of $f_j \forall j$. On the other hand, CMIM-2 [60] criterion corresponds exactly to the first term on the right hand side of eq. (34). These methods are able to detect pairs of relevant variables that act complementarily in predicting the class. In general CMIM-2 outperformed CMIM in experiments using artificial and benchmark datasets [60].

So far we have reviewed feature selection approaches that avoid estimating MI in high-dimensional spaces. Bonev *et al.* [9] proposed an extension of the MD criterion, called Max-min-Dependence (MmD), which is defined as follows:

$$J_{MmD}(f_i) = I(\{f_i, S\}; C) - I(\neg\{f_i, S\}; C). \quad (36)$$

The procedure starts with the empty set $S = \emptyset$ and sequentially generates S^{t+1} as:

$$S^{t+1} = S^t \cup \max_{f_i \in F \setminus S} (J_{MmD}(f_i)). \quad (37)$$

The MmD criterion is heuristic, and is not derived from a principled approach. However, Bonev *et al.* [9] were one of the first in selecting variables estimating MI in high-dimensional spaces [27], which allows using set of variables instead of individual variables. Chow and Huang [13] proposed combining a pruned Parzen window estimator with quadratic mutual information [47], using Renyi entropies, to estimate directly the MI between the feature subset S^t and the classes C , $I(S^t; C)$, in an effective and efficient way.

6 Open Problems

In this section we present some open problems and challenges in the field of feature selection, in particular from the point of view of information theoretic methods. Here can be found a non-exhaustive list of open problems or challenges.

- 1. Further developing a unifying framework for information theoretic feature selection.** As we reviewed in section 5 a unifying framework able to explain the advantages and limitations of successful heuristics has been proposed. This theoretical framework should be further developed in order to derive new efficient feature selection algorithms that include in their functional terms information related to the three types of features: relevant, redundant and complementary. Also a stronger connection between this framework and the Markov blanket is needed. Developing hybrid methods that combine maximal dependency with minimal conditional mutual information is another possibility.
- 2. Further improving the efficacy and efficiency of information theoretic feature selection methods in high-dimensional spaces.** The computational time depends on the search strategy and the evaluation criterion [24]. As we enter the era of Big Data, there is an urgent need for developing very fast feature selection methods able to work with millions of features and billions of samples. An important challenge is developing more efficient methods for estimating MI in high-dimensional spaces. Automatically determining the optimal size of the feature subset is also of interest, many feature selection methods do not have a stop criterion. Developing new search strategies that go beyond greedy optimization is another interesting possibility.
- 3. Further investigating the relationship between mutual information and Bayes error classification.** So far lower and upper bounds for error classification have been found for the case of one random variable and the target class. Extending these results to the case of mutual information between feature subsets and the target class is an interesting open problem.
- 4. Further investigating the effect of a finite sample over the statistical criteria employed and in MI estimation.** Guyon *et al.* [24] argued that feature subsets that are not sufficient may render better performance than sufficient feature subsets. For example, in the bio-informatics domain, it is common to have very large input dimensionality and small sample size [49].
- 5. Further developing a framework for studying the relation between feature selection and causal discovery.** Guyon *et al.* [25] investigated causal feature selection. The authors argued that the knowledge of causal relationships can benefit feature selection and viceversa. A challenge is to develop efficient Markov blanket induction algorithms for non-faithful distributions.

6. Developing new criteria of statistical dependence beyond correlation and MI. Seth and Principe [51] revised the postulates of measuring dependence according to Renyi, in the context of feature selection. An important topic is normalization, because a measure of dependence defined on different kinds of random variables should be comparable. There is no standard theory about MI normalization [19, 16]. Another problem is that estimators of measures of dependence should be good enough, even when using a few realizations, in the sense of following the desired properties of these measures. Seth and Principe [51] argued that this property is not satisfied by MI estimators, because they do not reach the maximum value under strict dependence, and are not invariant to one-to-one transformations.

7 Conclusions

We have presented a review of the state-of-the-art in information theoretic feature selection methods. We showed that modern feature selection methods must go beyond the concepts of relevance and redundancy to include complementarity (synergy). In particular, new feature selection methods that assess features in context are necessary. Recently, a unifying framework has been proposed, which is able to retrofit successful heuristic criteria. In this work, we have further developed this framework, presenting some new results and derivations. The unifying theoretical framework allows us to indicate the approximations made by each method, and therefore their limitations. A number of open problems in the field are suggested as challenges for the avid reader.

8 Acknowledgement

This work was funded by CONICYT-CHILE under grant FONDECYT 1110701.

References

1. Almuallim H, Dietterich TG (1991) Learning with many irrelevant features. In: Artificial Intelligence, Proceedings of the Ninth National Conference on, AAAI Press, pp 547–552
2. Almuallim H, Dietterich TG (1992) Efficient algorithms for identifying relevant features. In: Artificial Intelligence, Proceedings of the Ninth Canadian Conference on, Morgan Kaufmann, pp 38–45
3. Balagani K, Phoha V (2010) On the feature selection criterion based on an approximation of multidimensional mutual information. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 32(7):1342–1343
4. Battiti R (1994) Using mutual information for selecting features in supervised neural net learning. *Neural Networks, IEEE Transactions on* 5(4):537–550

5. Bell AJ (2003) The co-information lattice. *Analysis* pp 921–926
6. Bell DA, Wang H (2000) A formalism for relevance and its application in feature subset selection. *Machine Learning* 41(2):175–195
7. Bellman RE (1961) *Adaptive Control Processes: A Guided Tour*, 1st edn. Princeton University Press
8. Bins J, Draper B (2001) Feature selection from huge feature sets. In: *Computer Vision, 2001. Proceedings Eighth IEEE International Conference on*, vol 2, pp 159–165
9. Bonev B, Escolano F, Cazorla M (2008) Feature selection, mutual information, and the classification of high-dimensional patterns: Applications to image classification and microarray data analysis. *Pattern Analysis and Applications* 11(3-4):309–319
10. Brown G, Pocock A, Zhao MJ, Luján M (2012) Conditional likelihood maximisation: A unifying framework for information theoretic feature selection. *Journal of Machine Learning Research* 13:27–66
11. Brown LE, Tsamardinos I (2008) Markov blanket-based variable selection in feature space. Technical report dsl-08-01, Discovery Systems Laboratory, Vanderbilt University
12. Cheng H, Qin Z, Feng C, Wang Y, Li F (2011) Conditional mutual information-based feature selection analyzing for synergy and redundancy. *Electronics and Telecommunications Research Institute* 33(2):210–218
13. Chow T, Huang D (2005) Estimating optimal feature subsets using efficient estimation of high-dimensional mutual information. *Neural Networks, IEEE Transactions on* 16(1):213–224
14. Cover TM, Thomas JA (2006) *Elements of Information Theory*, 2nd edn. Wiley-Interscience
15. Davies S, Russell S (1994) N_p -completeness of searches for smallest possible feature sets. In: *Intelligent Relevance, Association for the Advancement of Artificial Intelligence Symposium on*, AAAI Press, pp 37–39
16. Duch W (2006) Filter methods. In: *Feature Extraction, Foundations and Applications, Studies in Fuzziness and Soft Computing*, vol 207, Springer Berlin Heidelberg, chap 3, pp 167–185
17. Duch W, Winiarski T, Biesiada J, Kachel A (2003) Feature selection and ranking filter. In: *Int. Conf. Artificial Neural Networks (ICANN) and Int. Conf. Neural Information Processing (ICONIP)*, pp 251–254
18. Estévez PA, Caballero R (1998) A niching genetic algorithm for selecting features for neural networks classifiers. In: *Perspectives in Neural Computation*, New York: Springer-Verlag, pp 311–316
19. Estévez PA, Tesmer M, Pérez CA, Zurada JM (2009) Normalized mutual information feature selection. *Neural Networks, IEEE Transactions on* 20(2):189–201
20. Feder M, Merhav N (1994) Relations between entropy and error probability. *Information Theory, IEEE Transactions on* 40(1):259–266
21. Fleuret F, Guyon I (2004) Fast binary feature selection with conditional mutual information. *Journal of Machine Learning Research* 5:1531–1555

22. Guo B, Nixon MS (2009) Gait feature subset selection by mutual information. *Systems, Man and Cybernetics, Part A, IEEE Transactions on* 39(1):36–46
23. Guyon I, Elisseeff A (2003) An introduction to variable and feature selection. *Journal of Machine Learning Research* 3:1157–1182
24. Guyon I, Elisseeff A (2006) An introduction to feature extraction. In: *Feature Extraction, Foundations and Applications, Studies in Fuzziness and Soft Computing*, vol 207, Springer Berlin Heidelberg, pp 1–25
25. Guyon I, Aliferis C, Elisseeff A (2008) Causal feature selection. In: Liu H, Motoda H (eds) *Computational Methods of Feature Selection*, Chapman & Hall/CRC, chap 4
26. Hellman M, Raviv J (1970) Probability of error, equivocation, and the chernoff bound. *Information Theory, IEEE Transactions on* 16(4):368–372
27. Hero A, Michel O (1999) Estimation of renyi information divergence via pruned minimal spanning trees. *Higher-Order Statistics Proceedings of the IEEE Signal Processing Workshop on* pp 264–268
28. Huang S (2003) Dimensionality reduction in automatic knowledge acquisition: a simple greedy search approach. *Knowledge and Data Engineering, IEEE Transactions on* 15(6):1364–1373
29. Jakulin A (2005) Learning based on attribute interactions. PhD thesis, University of Ljubljana, Slovenia
30. Jakulin A, Bratko I (2003) Quantifying and visualizing attribute interactions. *CoRR* cs.AI/0308002, URL <http://arxiv.org/abs/cs.AI/0308002>
31. Kohavi R, John GH (1997) Wrappers for feature subset selection. *Artificial Intelligence* 97(1-2):273–324
32. Kojadinovic I (2005) Relevance measures for subset variable selection in regression problems based on k-additive mutual information. *Computational Statistics and Data Analysis* 49(4):1205–1227
33. Koller D, Sahami M (1996) Toward optimal feature selection. Technical Report 1996-77, Stanford InfoLab
34. Kraskov A, Stögbauer H, Grassberger P (2004) Estimating mutual information. *Phys Rev E* 69:066,138
35. Kullback S (1997) *Information Theory and Statistics*. New York: Dover
36. Kullback S, Leibler RA (1951) On information and sufficiency. *Annals of Mathematical Statistics* 22:49–86
37. Kwak N, Choi CH (2002) Input feature selection for classification problems. *Neural Networks, IEEE Transactions on* 13(1):143–159
38. Lal KN, Chapelle O, Weston J, Elisseeff A (2006) Embedded methods. In: *Feature Extraction, Foundations and Applications, Studies in Fuzziness and Soft Computing*, vol 207, Springer Berlin Heidelberg, chap 5, pp 167–185
39. Lewis DD (1992) Feature selection and feature extraction for text categorization. In: *Proceedings of Speech and Natural Language Workshop*, Morgan Kaufmann, pp 212–217

40. Lin D, Tang X (2006) Conditional infomax learning: An integrated framework for feature extraction and fusion. In: Computer Vision - ECCV 2006, Lecture Notes in Computer Science, vol 3951, Springer Berlin Heidelberg, pp 68–82
41. Marill T, Green D (1963) On the effectiveness of receptors in recognition systems. *Information Theory, IEEE Transactions on* 9(1):11–17
42. McGill W (1954) Multivariate information transmission. *Psychometrika* 19(2):97–116
43. Meyer P, Schretter C, Bontempi G (2008) Information-theoretic feature selection in microarray data using variable complementarity. *Selected Topics in Signal Processing, IEEE Journal of* 2(3):261–274
44. Mo D, Huang SH (2011) Feature selection based on inference correlation. *Intelligent Data Analysis* 15(3):375–398
45. Mo D, Huang SH (2012) Fractal-based intrinsic dimension estimation and its application in dimensionality reduction. *Knowledge and Data Engineering, IEEE Transactions on* 24(1):59–71
46. Peng H, Long F, Ding C (2005) Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 27(8):1226–1238
47. Principe JC (2010) *Information Theoretic Learning: Renyi's Entropy and Kernel Perspectives*, 1st edn. Springer Publishing Company
48. Raudys S, Jain A (1991) Small sample size effects in statistical pattern recognition: recommendations for practitioners. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 13(3):252–264
49. Saeys Y, Inza I, Larranaga P (2007) A review of feature selection techniques in bioinformatics. *Bioinformatics* 23(19):2507–2517
50. Sebban M, Nock R (2002) A hybrid filter/wrapper approach of feature selection using information theory. *Pattern Recognition* 35(4):835–846
51. Seth S, Principe J (2010) Variable selection: A statistical dependence perspective. In: *Machine Learning and Applications (ICMLA)*, 2010 Ninth International Conference on, pp 931–936
52. Shannon CE (1948) A mathematical theory of communication. *Bell System Technical Journal* 27:379–423, 625–56
53. Somol P, Pudil P, Kittler J (2004) Fast branch & bound algorithms for optimal feature selection. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 26(7):900–912
54. Stearns SD (1976) On selecting features for pattern classifiers. In: *Pattern Recognition, Proceedings of the 3rd International Conference on*, Coronado, CA, pp 71–75
55. Tishby N, Pereira FC, Bialek W (1999) The information bottleneck method. *Proc 37th Annu Allerton Conf Communication, Control and Computing*
56. Torkkola K (2006) Information-theoretic methods. In: *Feature Extraction, Foundations and Applications, Studies in Fuzziness and Soft Computing*, vol 207, Springer Berlin Heidelberg, chap 6, pp 167–185

57. Trunk GV (1979) A problem of dimensionality: A simple example. *Pattern Analysis and Machine Intelligence*, IEEE Transactions on PAMI-1(3):306–307
58. Tsamardinos I, Aliferis CF (2003) Towards principled feature selection: Relevancy, filters and wrappers. In: *Artificial Intelligence and Statistics, Proceedings of the Ninth International Workshop on*, Morgan Kaufmann Publishers
59. Tsamardinos I, Aliferis CF, Statnikov E (2003) Algorithms for large scale markov blanket discovery. In: *The 16th International FLAIRS Conference*, St, AAAI Press, pp 376–380
60. Vergara JR, Estévez PA (2010) Cmim-2: An enhanced conditional mutual information maximization criterion for feature selection. *Journal of Applied Computer Science Methods* 2(1):5–20
61. Vidal-Naquet M, Ullman S (2003) Object recognition with informative features and linear classification. In: *Computer Vision, 2003. Proceedings Ninth IEEE International Conference on*, vol 1, pp 281–288
62. Watanabe S (1960) Information theoretical analysis of multivariate correlation. *IBM Journal of Research and Development* 4(1):66–82
63. Webb AR (2002) *Statistical Pattern Recognition*, 2nd edn. John Wiley & Song, Ltd.
64. Weston J, Mukherjee S, Chapelle O, Pontil M, Poggio T, Vapnik V (2000) Feature selection for svms. In: *Advances in Neural Information Processing Systems 13*, MIT Press, pp 668–674
65. Whitney A (1971) A direct method of nonparametric measurement selection. *Computers*, IEEE Transactions on C-20(9):1100–1103
66. Yang HH, Moody J (1999) Feature selection based on joint mutual information. In: *Advances in Intelligent Data Analysis, Proceedings of International ICSC Symposium on*, pp 22–25
67. Yu L, Liu H (2004) Efficient feature selection via analysis of relevance and redundancy. *Journal of Machine Learning Research* 5:1205–1224
68. Zhao Z, Liu H (2009) Searching for interacting features in subset selection. *Intelligent Data Analysis* 13(2):207–228